# An Audit of Misinformation Filter Bubbles on YouTube: Bubble Bursting and Recent Behavior Changes

Matus Tomlein, Branislav Pecher, Jakub Simko, Ivan Srba, Robert Moro, Elena Stefancova, Michal Kompan, Andrea Hrckova, Juraj Podrouzek, Maria Bielikova

Stuart Heeb

# Have you ever noticed…

# Let's take a trip down the YouTube rabbit hole

YouTube ᶜᴴ

andrew tate

Anmelden

Startseite

Shorts

Abos

Mediathek

Verlauf

Melde dich an, um Videos mit "Mag ich" zu bewerten, zu kommentieren und um Kanäle zu abonnieren.

Anmelden

Entdecken

Trends

Musik

Filme & Serien

Gaming

Sport

Kanäle finden

Mehr von YouTube

YouTube Premium

YouTube Music

YouTube Kids

Filter



**Andrew Tate First Words After Release Gets REJECTED**

183.607 Aufrufe • vor 22 Stunden

Fesify ✓

Andrew Tate First Words After Release Gets REJECTED ▷ Click Here To Subscribe ▷ http://goo.gl/Q6lJeG Fesify is a media ...

Neu

Andrew Tate NOT Getting Released Now (Court Update)

Fesify ✓

Katja Krasavice EXPOSED | Andrew Tate ist TODKRANK?!

Andrew Tate is having a bad time right now.

**Andrew Tate First Words After Release Gets REJECTED**

**Fesify** ✓
328.000 Abonnenten

Abonnieren

👍 4682  👎  ➤ Teilen  Speichern  •••

183.607 Aufrufe  vor 22 Stunden
Andrew Tate First Words After Release Gets REJECTED

▶ Click Here To Subscribe ▶ http://goo.gl/Q6lJeG Mehr ansehen

1.799 Kommentare    Sortieren nach

Kommentar hinzufügen...

andrew tate

Anmelden

**Andrew Tate RUSHED To Hospital With Emergency…**
Fesify ✓
35.545 Aufrufe · vor 5 Tagen
Neu

6

# Recommendations

"Our recommendation system is built on the simple principle of helping people find the videos they **want to watch** and that will **give them value**"

**Exploration vs. exploitation**

# Filter bubbles

**DeepMind** ✅
@DeepMind

Feedback loops in recommendation systems can give rise to "echo chambers" and "filter bubbles" which can narrow a user's content exposure, and ultimately shift their world view.

5:06 PM · Mar 1, 2019

https://twitter.com/DeepMind/status/1101514121563041792?s=20

generated by DALL-E

Is bubble bursting possible?

# Motivation

- Need for independent oversight of personalization behavior



generated by DALL-E

# Reference study

- Hussein et al. (2020): **Measuring Misinformation in Video Search Platforms: An Audit Study on YouTube** [1], experiment conducted in mid 2019

- This study's experiment was conducted in March 2021

# Reference study

- Filter bubbles are easily and quickly created

- "YouTube still has a long way to go to mitigate misinformation on its platform" [1]

# Audits

**Crowdsourcing**

- using real user data
- uncontrolled environment
- hard to make comparisons

**Sockpuppeting**

# Audits

## Crowdsourcing

- using real user data
- uncontrolled environment
- hard to make comparisons

## Sockpuppeting

- using non-human bots
- selection of appropriate seed data

# Agents



N. Virginia

06.06.1990

„rather not say"

HOME PAGE

SEARCH RESULTS

UP-NEXT
RECOMMENDATIONS

# Agents

- Watches videos for ≤ 30 mins

- Does **not**

  - Like

  - Subscribe

  - Comment

  - **Act human!**

# Human factors

- Selective exposure

- Confirmation bias

- Dunning-Kruger effect

# Human factors

- **Selective exposure**

- Confirmation bias

- Dunning-Kruger effect

generated by DALL-E

# Human factors

- Selective exposure

- **Confirmation bias**

- Dunning-Kruger effect

**Objective facts**

**What confirms your beliefs**

**What you see**

# Human factors

- Selective exposure

- Confirmation bias

- **Dunning-Kruger effect**

**"The less you know, the more confident you are"**

# Experiment

Agent Initialization → Promoting → Debunking → Tear-down

# Experiment



**Agent Initialization** → **Promoting** → **Debunking** → **Tear-down**

**Phase 0**

- Most popular promoting/debunking videos (seed data)

- Search queries (e.g. "9/11 conspiracy")

- Wait time between each query

# Experiment

| Agent Initialization | Promoting | Debunking | Tear-down |

**Phase 1**

- Create the filter bubble

# Experiment

- Burst the filter bubble

# Experiment

| Agent Initialization | Promoting | Debunking | Tear-down |
|---|---|---|---|

**Phase 3**

- Clear YouTube history

# Topics

9/11

Chem-trails

Anti-vaccination

Moon landing

Flat earth

# Metrics

- Score $x_i$ of a single video

$$x_i = 1 \Leftrightarrow \quad \text{promoting}$$
$$x_i = 0 \Leftrightarrow \quad \text{neutral}$$
$$x_i = -1 \Leftrightarrow \quad \text{debunking}$$

# Metrics

- **N**ormalized **S**core

$$\text{NS} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

**For recommendations**

# NS



$$\mathrm{NS} = \frac{1}{n}\sum_{i=1}^{n} x_i$$

$$= \frac{1}{n}(x_1 + x_2 + x_3 + x_4 + x_5$$

$$+ x_6 + x_7 + x_8 + x_9 + x_{10})$$

$$= \frac{1}{10}(1 + 0 - 1 + 1 + 0$$

$$+ 1 + 0 + 1 - 1 + 1) = 0.3$$

# NS



$$NS = \frac{1}{n}\sum_{i=1}^{n} x_i$$

neutral

debunking

$-1$

$0.3$

$+1$

promoting

# Metrics

- **Se**arch **R**esult **P**age **M**isinformation **S**core

$$\text{SERP-MS} = \frac{1}{\frac{n \cdot (n+1)}{2}} \sum_{i=1}^{n} x_i \cdot (n - r_i + 1)$$

rank of video $i$

**For search results**

# SERP-MS



$$\text{SERP-MS} = \frac{1}{\frac{n \cdot (n+1)}{2}} \sum_{i=1}^{n} x_i \cdot (n - r_i + 1)$$

$$= \frac{1}{\frac{n \cdot (n+1)}{2}} \big( \quad x_1 \cdot (n - r_1 + 1)$$
$$+ x_2 \cdot (n - r_2 + 1)$$
$$+ x_3 \cdot (n - r_3 + 1)$$
$$+ x_4 \cdot (n - r_4 + 1)$$
$$+ x_5 \cdot (n - r_5 + 1) \big)$$

# SERP-MS



$$\text{SERP-MS} = \frac{1}{\frac{n \cdot (n+1)}{2}} \sum_{i=1}^{n} x_i \cdot (n - r_i + 1)$$

$$= \frac{1}{15} \left( \begin{array}{l} 1 \cdot (5 - 1 + 1) \\ -1 \cdot (5 - 2 + 1) \\ +0 \cdot (5 - 3 + 1) \\ -1 \cdot (5 - 4 + 1) \\ -1 \cdot (5 - 5 + 1) \end{array} \right)$$

$$= \frac{1}{15}(5 - 4 + 0 - 2 - 1) = -0.133$$

# SERP-MS



neutral

debunking $\qquad$ promoting

$-1 \qquad -1.33 \qquad +1$

# Hypotheses

"How has YouTube's personalization behavior changed with regards to misinformation videos since the reference study?" [2]

# Results

- Comparison with reference study (expecting **better**)

| Topic | Search results score | Recommendation score |
|---|---|---|
| 9/11 | n.s.d. | n.s.d. |
| Chemtrails | n.s.d. | n.s.d. |
| Flat earth | n.s.d. | n.s.d. |
| Moon landing | n.s.d. | **better** |
| Anti-vaccination | **worse** — Less debunking videos — | **worse** |

n.s.d. = not statistically significantly different

# Hypotheses

"How does the effect of misinformation filter bubbles change,

when debunking videos are watched?" [2]

# Results

- Bubble creating behavior (expecting **worse**)

| Topic | Search results score | Recommendation score |
|---|---|---|
| 9/11 | n.s.d. | **worse** |
| Chemtrails | n.s.d. | n.s.d. |
| Flat earth | **better** ── Promoting videos disappear in some queries | n.s.d. |
| Moon landing | n.s.d. | n.s.d. |
| Anti-vaccination | n.s.d. | **worse** |

n.s.d. = not statistically significantly different

# Results



- Bubble bursting behavior (expecting **better**)

| Topic | Search results score | Recommendation score |
|---|---|---|
| 9/11 | n.s.d. | **better** |
| Chemtrails | n.s.d. | **better** |
| Flat earth | n.s.d. | **better** |
| Moon landing | n.s.d. | n.s.d. |
| Anti-vaccination | **better** | **better** |

n.s.d. = not statistically significantly different

# Results

| | Agent Initialization | Promoting | Debunking | Tear-down |
|---|---|---|---|---|

- Comparison to baseline (expecting **better**)

| Topic | Search results score | Recommendation score |
|---|---|---|
| 9/11 | n.s.d. | n.s.d. |
| Chemtrails | **better** | **better** |
| Flat earth | **better** | **better** |
| Moon landing | **better** | n.s.d. |
| Anti-vaccination | **better** | **better** |

n.s.d. = not statistically significantly different

# Outlook

- Srba et al. (2023): **Auditing YouTube's Recommendation Algorithm for Misinformation Filter Bubbles** [3], continuation of this paper

# Conclusion

- YouTube seems to have not fulfilled its pledges

- Bubble bursting is possible, but there are differences between topics

# Discussion

- What *is* misinformation?

- How much should YouTube intervene in this matter?

- How strongly should recommendations adhere to human tendencies?

- Does YouTube treat misinformation topics differently?

- Study annotation score vs. YouTube's "internal scoring"

# References

[1]  Reference study: https://dl.acm.org/doi/10.1145/3392854 (2020)

[2]  Paper: https://dl.acm.org/doi/pdf/10.1145/3460231.3474241 (2021)

[3]  Continuation of paper: https://arxiv.org/abs/2210.10085 (2023)

# Sources

- Downward stair case: DALL-E (http://labs.openai.com), „a downward spiral into a dark dimension, digital art" (accessed March 14, 2023)

- YouTube Recommendation System: https://blog.youtube/inside-youtube/on-youtubes-recommendation-system/ (accessed March 13, 2023)

- Screenshot of tweet: https://twitter.com/DeepMind/status/1101514121563041792?s=20 (accessed March 14, 2023)

- Bubble: DALL-E (http://labs.openai.com), „a soap bubble, that is also a portal to a dark dimension, digital art" (accessed March 13, 2023)

- Bubble bursting image: https://www.inc.com/partners-in-leadership/4-strategies-to-burst-your-filter-bubble-and-influence-others.html (accessed March 13, 2023)

- Computer surveillance: DALL-E (http://labs.openai.com), „dystopian 1984 themed image of computer activity being surveilled, digital art" (accessed March 15, 2023)

- YouTube's pledges (October 15, 2020): https://blog.youtube/news-and-events/harmful-conspiracy-theories-youtube/ (accessed March 13, 2023)

- Gender-neutral icon: https://thenounproject.com/icon/gender-neutral-147092/ (accessed March 13, 2023)

- Magnifying glass: DALL-E (http://labs.openai.com), „magnifying glass which puts only a part of a document in focus" (accessed March 19, 2023)